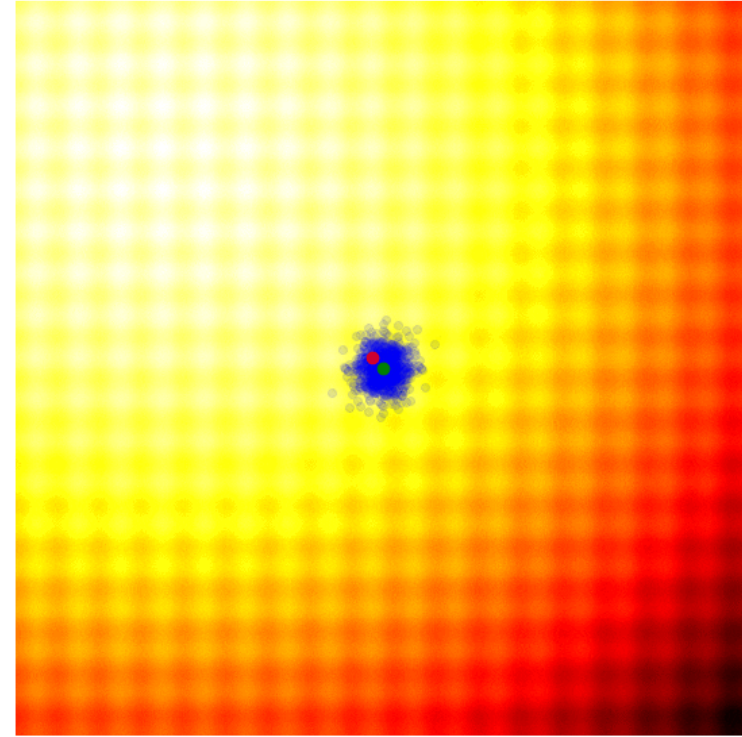
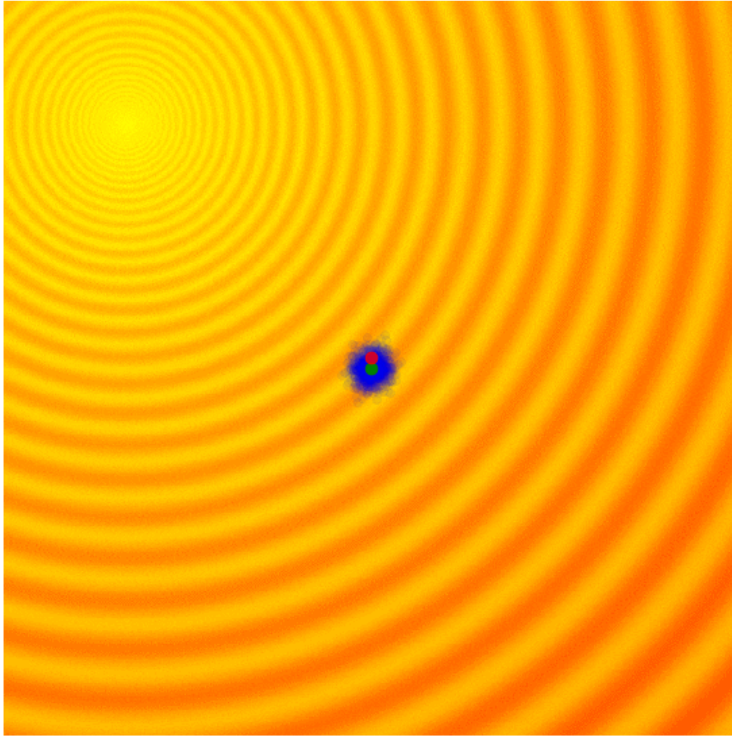


# Introduction to Reinforcement Learning

Lecture 5. Population-based Method

Sungjoon Choi, Korea University

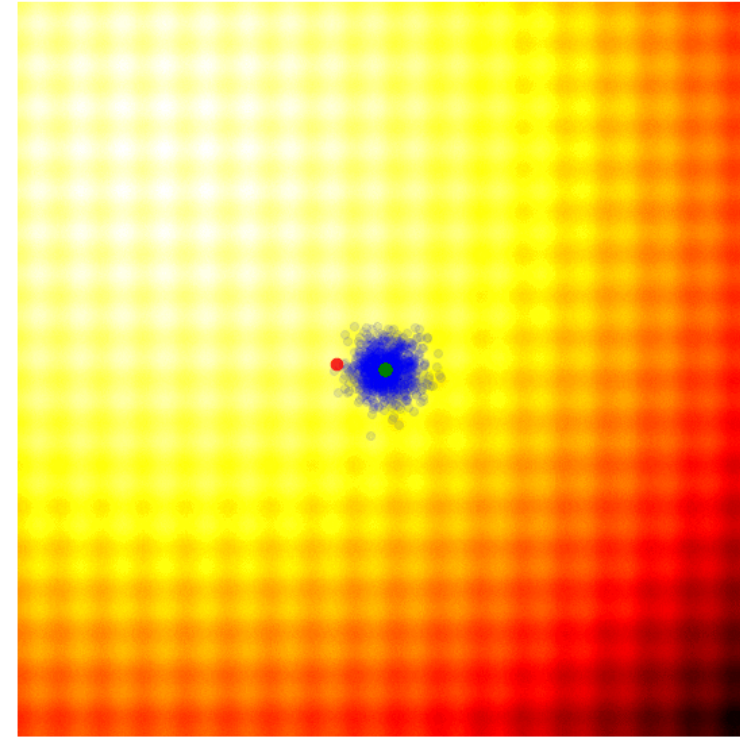
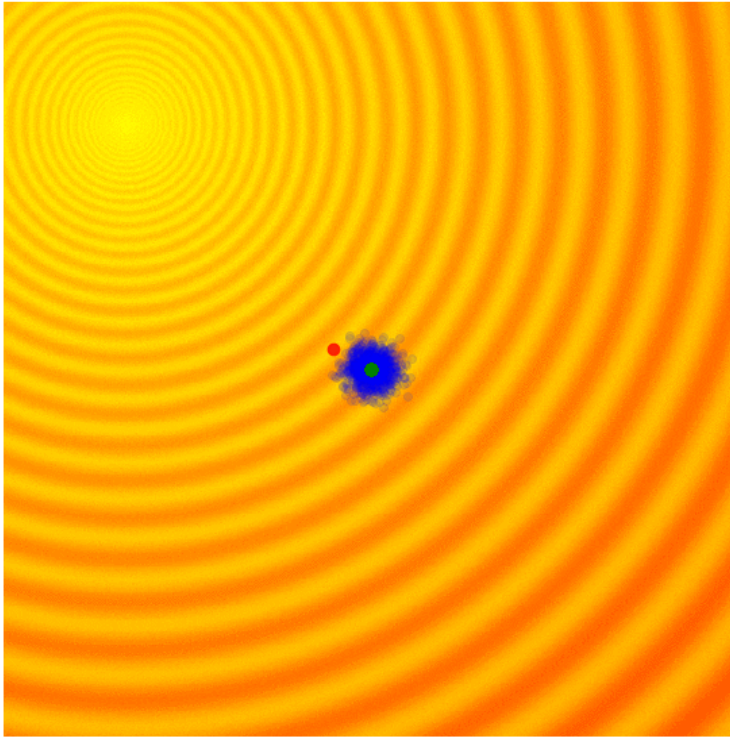
# Population-based Search



- The **green dot** indicates the mean of the distribution at each generation, the **blue dots** are the sampled solutions, and the **red dot** is the best solution so far.

# Cross Entropy Method (CEM)

# Genetic Algorithm



- The **green dot** indicates the **elite population** from the previous generation, the **blue dots** are the offsprings to form the set of candidate solutions, and the **red dot** is the best solution so far.

# Cross Entropy Method (CEM)

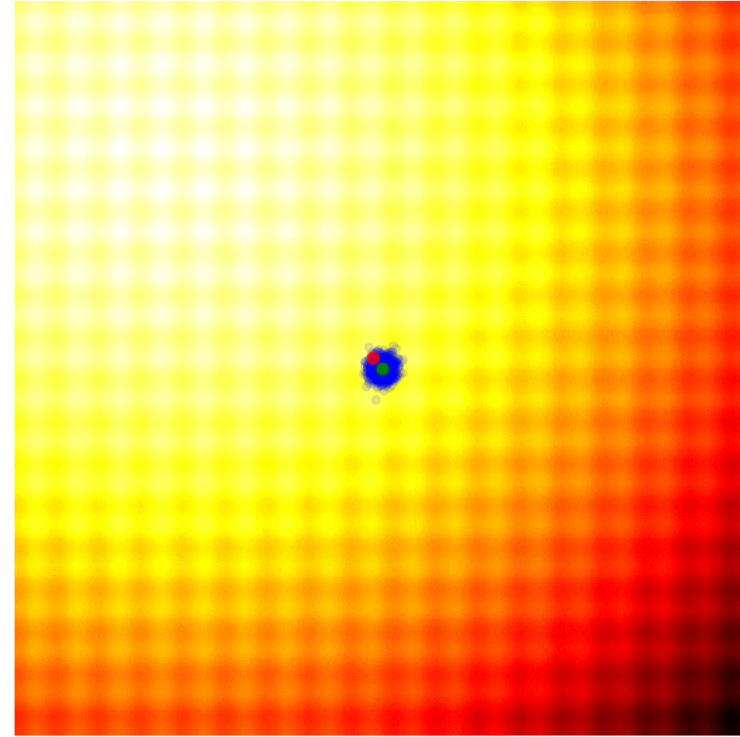
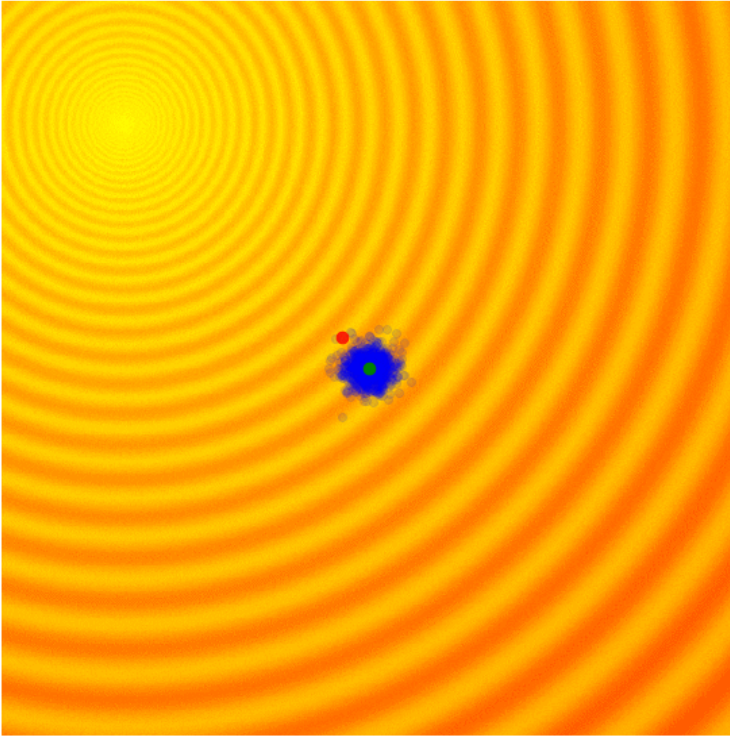


- Initialize  $\mu \in \mathbb{R}^d$  and  $\sigma \in \mathbb{R}^d$
- **for** iteration = 1, 2, ...
  - Collect  $n$  samples of  $\theta_i \sim \mathcal{N}(\mu, \text{diag}(\sigma))$
  - Perform a noisy evaluation  $R_i \sim \theta_i$
  - Select the top  $p\%$  of samples (aka the **elite set**)
  - Fit a Gaussian distribution (with diagonal covariance) to the **elite set**, obtaining a new  $\mu$  and  $\sigma$
- **end for**
- Return the final  $\mu$

# CMA-ES

"The CMA Evolution Strategy: A Tutorial," 2016

# CMA-ES



- The **green dot** indicates the **elite population** from the previous generation, the **blue dots** are the offsprings to form the set of candidate solutions, and the **red dot** is the best solution so far.
- CMA-ES can **adaptively** increase or decrease the search space!

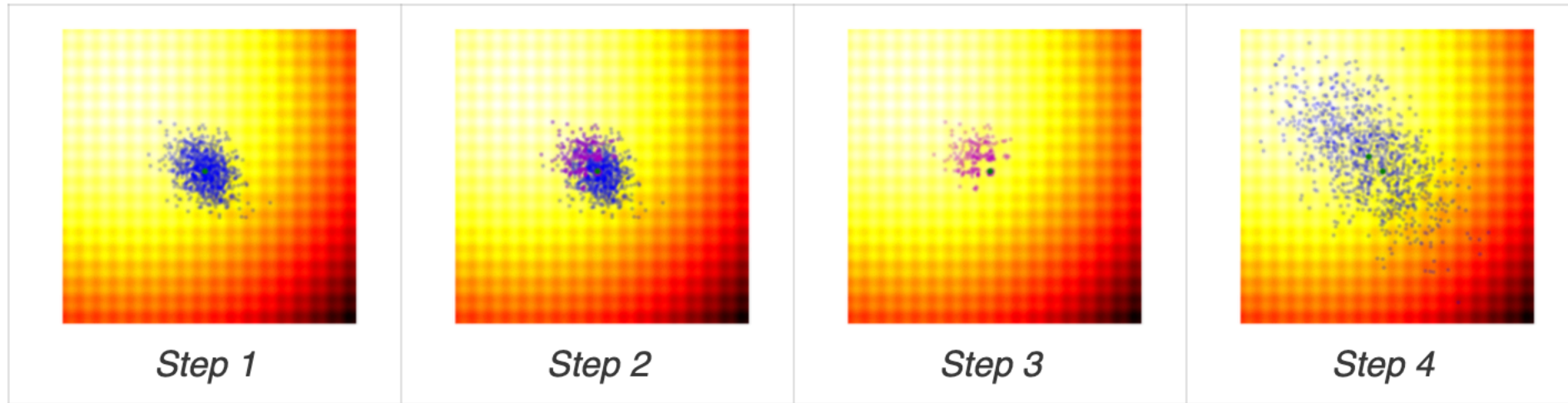
# CMA-ES



- Initialize  $\mu \in \mathbb{R}^d$  and  $\sigma \in \mathbb{R}^d$
- **for** iteration  $g = 1, 2, \dots$ 
  - Collect  $n$  samples of  $\theta_i \sim \mathcal{N}(\mu, \text{diag}(\sigma))$
  - Perform a noisy evaluation  $R_i \sim \theta_i$
  - Select the top  $p$  % of samples (aka the **elite set**)
  - Fit  $\mu^{(g)}$  to the **elite set** by  $\mu^{(g)} = \frac{1}{N_{\text{best}}} \sum_i^{N_{\text{best}}} \theta_i$
  - The trick is to update the covariance using the previous mean  $\mu^{(g-1)}$ .
    - For 1-dimensional case,  $\sigma_x^{2,(g+1)} = \frac{1}{N_{\text{best}}} \sum_{i=1}^{N_{\text{best}}} (x_i - \mu_x^{(g)})^2$
- **end for**
- Return the final  $\mu$



# The Effect of CMA



# Augmented Random Search (ARS)

"Simple random search provides a competitive approach to reinforcement learning," 2018

# Augmented Random Search



- Augmented Random Search (ARS) is a population-based, derivative-free reinforcement learning method.
  - Derivative-free methods such as a **cross-entropy method** (CEM) or **covariance matrix adaptation** (CMA) treat the return as a black box function to be optimized in terms of the policy parameters.
  - From the current policy parameter  $\theta^t$ , multiple parameters  $\{\tilde{\theta}_i^t\}_{i=1}^N$  are sampled and the returns  $\{\eta(\pi_{\tilde{\theta}_i^t})\}_{i=1}^N$  of the corresponding policies  $\{\pi_{\tilde{\theta}_i^t}\}_{i=1}^N$  evaluated by rollouts.
- ARS utilizes a simple linear policy and searches over the space of matrices.

# Basic Random Search



---

## Algorithm 1 Basic Random Search (BRS)

---

- 1: **Hyperparameters:** step-size  $\alpha$ , number of directions sampled per iteration  $N$ , standard deviation of the exploration noise  $\nu$
- 2: **Initialize:**  $\theta_0 = \mathbf{0}$ , and  $j = 0$ .
- 3: **while** ending condition not satisfied **do**
- 4:   Sample  $\delta_1, \delta_2, \dots, \delta_N$  of the same size as  $\theta_j$ , with i.i.d. standard normal entries.
- 5:   Collect  $2N$  rollouts of horizon  $H$  and their corresponding rewards using the policies

$$\pi_{j,k,+}(x) = \pi_{\theta_j + \nu \delta_k}(x) \quad \text{and} \quad \pi_{j,k,-}(x) = \pi_{\theta_j - \nu \delta_k}(x),$$

Positively perturbed policy                      Negatively perturbed policy

with  $k \in \{1, 2, \dots, N\}$ .

- 6:   Make the update step:

$$\theta_{j+1} = \theta_j + \underbrace{\frac{\alpha}{N}}_{\text{Stepsize}} \sum_{k=1}^N \underbrace{[r(\pi_{j,k,+}) - r(\pi_{j,k,-})]}_{\text{Score of the perturbation}} \underbrace{\delta_k}_{\text{Perturbation}}.$$

- 7:    $j \leftarrow j + 1$ .
  - 8: **end while**
-

# Augmented Random Search



- Drawbacks of Basic Random Search
  - Input normalization is often necessary for high-dimensional inputs.
  - Not all perturbations (directions) are useful.

# Augmented Random Search



**Algorithm 2** Augmented Random Search (ARS): four versions **V1**, **V1-t**, **V2** and **V2-t**

- 1: **Hyperparameters:** step-size  $\alpha$ , number of directions sampled per iteration  $N$ , standard deviation of the exploration noise  $\nu$ , number of top-performing directions to use  $b$  ( $b < N$  is allowed only for **V1-t** and **V2-t**) Use top  $b$  search directions.
- 2: **Initialize:**  $M_0 = \mathbf{0} \in \mathbb{R}^{p \times n}$ ,  $\mu_0 = \mathbf{0} \in \mathbb{R}^n$ , and  $\Sigma_0 = \mathbf{I}_n \in \mathbb{R}^{n \times n}$ ,  $j = 0$ .
- 3: **while** ending condition not satisfied **do**
- 4:   Sample  $\delta_1, \delta_2, \dots, \delta_N$  in  $\mathbb{R}^{p \times n}$  with i.i.d. standard normal entries.
- 5:   Collect  $2N$  rollouts of horizon  $H$  and their corresponding rewards using the  $2N$  policies

$$\begin{aligned} \mathbf{V1}: \quad & \begin{cases} \pi_{j,k,+}(x) = (M_j + \nu\delta_k)x \\ \pi_{j,k,-}(x) = (M_j - \nu\delta_k)x \end{cases} \\ \mathbf{V2}: \quad & \begin{cases} \pi_{j,k,+}(x) = (M_j + \nu\delta_k) \text{diag}(\Sigma_j)^{-1/2}(x - \mu_j) \\ \pi_{j,k,-}(x) = (M_j - \nu\delta_k) \text{diag}(\Sigma_j)^{-1/2}(x - \mu_j) \end{cases} \end{aligned} \quad \text{Input normalization}$$

for  $k \in \{1, 2, \dots, N\}$ .

- 6:   Sort the directions  $\delta_k$  by  $\max\{r(\pi_{j,k,+}), r(\pi_{j,k,-})\}$  denote by  $\delta_{(k)}$  the  $k$ -th largest direction, and by  $\pi_{j,(k),+}$  and  $\pi_{j,(k),-}$  the corresponding policies.
- 7:   Make the update step:

$$M_{j+1} = M_j + \frac{\alpha}{b\sigma_R} \sum_{k=1}^b [r(\pi_{j,(k),+}) - r(\pi_{j,(k),-})] \delta_{(k)},$$

where  $\sigma_R$  is the standard deviation of the  $2b$  rewards used in the update step.

- 8:   **V2** : Set  $\mu_{j+1}$ ,  $\Sigma_{j+1}$  to be the mean and covariance of the  $2NH(j+1)$  states encountered from the start of training.<sup>2</sup>
- 9:    $j \leftarrow j + 1$
- 10: **end while**



# Thank You



ROBOT INTELLIGENCE LAB